

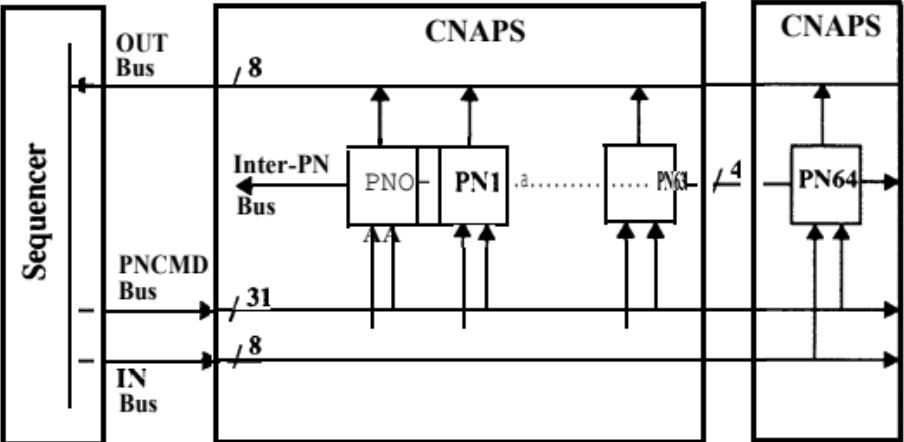
# Exhibit 4

## Exhibit 3 - CNAPS



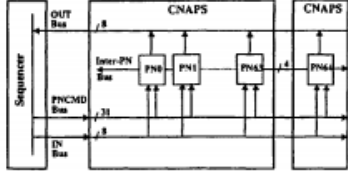
'156 Patent

Claim Limitation (Claim 7)	Exemplary Disclosure
<p>[156a] A device comprising:</p>	<p>Connected Network of Adaptive Processors (“CNAPS”) discloses a device. Specifically, the CNAPS is computer architecture consisting of an array of processors controlled by a sequencer. <i>See, e.g.:</i></p> <p>“The CNAPS architecture consists of an array of processors controlled by a sequencer, both implemented as a chip set developed by Adaptive Solutions. Inc. The sequencer is a one-chip device called the CNAPS Sequencer Chip (CSC). The processor array is also a one-chip device, available with either 64 or 16 processors per chip (the CNAPS- 1064 or CNAPS- 10 16). The CSC can control up to eight 1064s or 1016s, which act like <i>one large device</i>.” Dan Hammerstrom, AN INTRODUCTION TO NEURAL AND ELECTRONIC NETWORKS 336 (2d ed. 1995) (emphasis added).</p> <p>“The digital CNAPS neurocomputer (Adaptive Solutions, Inc.) has been designed as a general purpose parallel computer for artificial neural networks and image processing tasks and was introduced by Hammerstrom in 1990. CNAPS (Connected Network of Adaptive Processors) consists of a sequencer and a linear array of maximal processor nodes and is a SIMD computer (single instruction stream – multiple data stream). A single instruction bus and three data busses connect the sequencer and the processor nodes: a broadcast input bus (8 bit), a shared output bus (8 bit), and a circular inter-PN bus (2 bit) connecting neighbouring processors (figure 3 left). The sequencer broadcasts an instruction to the processor nodes single instruction stream which synchronously execute it but process the data from their local memories multiple data stream.” Peter Paschke &amp; Ralf Möller, <i>Simulation of Sparse Random Networks on a CNAPS SIMD Neurocomputer</i>, Proc. 15th IMACS World Congress, Berlin, August 1997, at 766.</p>

## Exhibit 3 - CNAPS

Claim Limitation (Claim 7)	Exemplary Disclosure
	 <p data-bbox="856 711 1753 792">FIGURE 1 The basic CNAPS architecture. CNAPS is a single instruction, multiple data (SIMD) architecture that uses broadcast input, one-dimensional interprocessor communication, and a single shared output bus.</p> <p data-bbox="693 816 1879 881">Dan Hammerstrom, AN INTRODUCTION TO NEURAL AND ELECTRONIC NETWORKS 338 (2d ed. 1995).</p> <p data-bbox="693 925 1906 1393">“The Adaptive Solutions CNAPS architecture is embodied in a single chip digital neurocomputer with 64 processors running at 25 megahertz. All processors receive the same instruction which they conditionally execute. Multiplication and addition are performed in parallel allowing 1.6 billion inner product steps per second per chip. Each processor has a 32-bit adder, 9-bit by 16-bit multiplier (16 by 16 in two clock cycles), shifter, logic unit, 32 16-bit registers, and 4096 bytes of local memory. Input and output are accomplished over 8-bit input and output buses common to all processors. The output bus is tied to the input bus so that output of one processor can be broadcast to all others. When multiple chips are used, they appear to the user as one chip with more processors. Special circuits support finding the maximum of values held in each processor and conserving weight space for sparsely connected networks. An accompanying sequencer chip controls instruction flow, input and output.” Hal McCartor, <i>Back Propagation Implementation on the Adaptive Solutions CNAPS Neurocomputer Chip</i>, ADVANCES IN NEURAL INFORMATION PROCESSING SYSTEMS 3, 1990, at</p>

## Exhibit 3 - CNAPS

Claim Limitation (Claim 7)	Exemplary Disclosure
	<p>1029.</p> <div style="display: flex; justify-content: space-around; align-items: flex-start;"> <div style="text-align: center;">  <p><b>Adaptive Solutions</b></p> <p><b>THE ARCHITECTURE</b></p> <ul style="list-style-type: none"> <li>• A linear array of PNs</li> <li>• SIMD control, each PN executes the same instruction each clock</li> <li>• External sequencer with writeable program store</li> <li>• Broadcast instruction bus (31 bits), PNCMD</li> <li>• Broadcast input data bus (8 bits), IN bus</li> <li>• Single, arbitrated output bus (8 bits), OUT bus</li> </ul> </div> <div style="text-align: center;">  <p><b>Adaptive Solutions</b></p> <p><b>Basic CNAPS Array Layout:</b></p>  </div> </div> <p>Dan Hammerstrom &amp; Gary Tahara, <i>CNAPS (Connected Network of Adaptive Processors)</i>, Presentation at Hot Chips Conference, July 8, 1991, <a href="https://www.hotchips.org/wp-content/uploads/hc_archives/hc03/3_Tue/HC3.S7/HC3.7.2.pdf">https://www.hotchips.org/wp-content/uploads/hc_archives/hc03/3_Tue/HC3.S7/HC3.7.2.pdf</a>.</p>
<p>[156b] at least one first low precision high dynamic range (LPHDR) execution unit adapted to execute a first operation on a first input signal representing a first numerical value to produce a first output signal representing a second numerical value,</p>	<p>CNAPS discloses at least one first low precision high dynamic range (LPHDR) execution unit adapted to execute a first operation on a first input signal representing a first numerical value to produce a first output signal representing a second numerical value. <i>See, e.g.:</i></p> <p>“The Adaptive Solutions CNAPS architecture is embodied in a single chip digital neurocomputer with 64 processors running at 25 megahertz. All processors receive the same instruction which they conditionally execute. Multiplication and addition are performed in parallel allowing 1.6 billion inner product steps per second per chip. Each processor has a 32-bit adder, 9-bit by 16-bit multiplier (16 by 16 in two clock cycles), shifter, logic unit, 32 16-bit registers, and 4096 bytes of local memory. Input and output are accomplished over 8-bit input and output buses common to all processors.” Hal McCartor, <i>Back Propagation Implementation on the Adaptive Solutions CNAPS Neurocomputer Chip</i>, <i>ADVANCES IN NEURAL INFORMATION PROCESSING SYSTEMS</i> 3, 1990, at 1029.</p> <p>“CNAPS is a single instruction, multiple data stream (SIMD) architecture. A SIMD computer</p>

**Exhibit 3 - CNAPS**

Claim Limitation (Claim 7)	Exemplary Disclosure
	<p>has one instruction sequencing/control unit and many processor nodes (PNs). In CNAPS, the PNs are connected in a one-dimensional array (Figure 1) in which each PN can ‘talk’ only to its right or left neighbors. The sequencer broadcasts each instruction plus input data to all PNs, which execute the same instruction at each clock. The PNs transmit output data to the sequencer, with several arbitration modes controlling access to the output bus.</p> <p>As Figure 2 suggests, each PN has a local memory, a multiplier, an adder/subtractor, a shifter/logic unit, a register tile, and a memory addressing unit. The entire PN uses fixed-point, two’s complement arithmetic, and the precision is 16 bits, with some exceptions. The PN memory can handle 8- or 16-bit reads or writes. The multiplier produces a 24-bit output; an 8 X 16 or 8 X 8 multiply takes one clock, and a 16 X 16 multiply takes two clocks. The adder can switch between 16- or 32-bit modes. The input and output buses are 8 bits wide, and a 16-bit word can be assembled (or disassembled) from two bytes in two clocks.” Dan Hammerstrom, AN INTRODUCTION TO NEURAL AND ELECTRONIC NETWORKS 337 (2d ed. 1995).</p>

## Exhibit 3 - CNAPS

Claim Limitation (Claim 7)	Exemplary Disclosure
	<p>FIGURE 2 The internal structure of a CNAPS processor node (PN). Each PN has its own storage and arithmetic capabilities. Storage consists of 4096 bytes. Arithmetic operations include multiply, accumulate, logic, and shift. All units are interconnected by two 16-bit buses.</p> <p>Dan Hammerstrom, AN INTRODUCTION TO NEURAL AND ELECTRONIC NETWORKS 338 (2d ed. 1995).</p> <p>“The X1 system consists of a linear array of PNs. Each PN is a simple arithmetic processor with its own local memory. The array is sequenced by a single controller, thus each PN executes the same instruction at each clock, using SIMD processing. Each PN is connected to three global buses: InBus—the data input bus, PnCmd—the command bus, which indicates what operations a PN performs each clock, and OutBus - the output bus.” Dan Hammerstrom, <i>A</i></p>

## Exhibit 3 - CNAPS

Claim Limitation (Claim 7)	Exemplary Disclosure
	<p><i>VLSI architecture for high-performance, low-cost, on-chip learning</i>, 1990 IJCNN International Joint Conference on Neural Networks, 1990, vol.2, at 539.</p> <p>“Each Processor Node has internal units connected via control signals and several data buses. See Figure 5. The InBus and PnCmd bus enter the PN through the Input Unit. In general the PN is horizontally microcoded:</p> <ol style="list-style-type: none"> <li>1. <i>Input Unit</i>: The Input Unit decodes the PnCmd bus and routes it to the other units, and receives an 8 bit value from the InBus. It also contains a flag to allow conditional instruction execution of each PN. Although the input and output buses are 8 bits, the internal buses are 16 bits. The Input Unit can assemble two 8 bit quantities into a 16 bit value when needed.</li> <li>2. <i>Logic-shifter</i>: The Logic-shifter contains both a shifter and a logic operation unit. Both operate on 16 bit quantities. These allow computed quantities to be manipulated by shifting and bit masking. There is also a 1’s counter that sums the bits in the byte and whose output is used with 1 bit weights (8 per byte) to accumulate the AND of an 8 bit input and an 8 bit weight. This permits binary inputs and binary weights (8 per byte) - 8 are computed each clock.</li> <li>3. <i>Registerfile</i>: The register file contains 32 16-bit registers for intermediate storage of constants such as learning rates and the PN ID.</li> <li>4. <i>Multiplier</i>: The Multiplier unit contains a 9x16 bit, two’s complement multiplier. The multiplier produces a 24 bit output. Using mode bits, the personality of the multiplier can be modified to multiply a positive 8 bit number by a signed 16 bit number, a signed 8 bit number by a signed 16 bit number, or a signed 16 bit number by a signed 16 bit number (this operation takes two clocks). There is a direct path from the multiplier to the adder that allows for the PN to be operated in Vector Mode where a simultaneous multiply and accumulate can occur each clock</li> <li>5. <i>Adder</i>: The adder/subtractor takes 32 bit inputs and produces a two’s complement 32 bit result. Adder overflow causes saturation to the largest positive or negative number, depending on the sign of the final result.</li> <li>6. <i>Weight Address Generation</i>: The separate multiplier-to-adder bus allows the PN to be operated in a vector mode, which requires a regular stream of addresses to be generated for the weight memory. This unit contains its own adder for adding the contents of the stride</li> </ol>

## Exhibit 3 - CNAPS

Claim Limitation (Claim 7)	Exemplary Disclosure
	<p>register to the current weight index. As a result there is arbitrary striding through memory for those programs that have more complex data structures.</p> <p>7. <i>Weight Memory</i>: Memory can be accessed in either 8 bit or 16 bit mode. There is also a hardware system, called the Virtual Zeros mechanism<sup>4</sup>, which is used for the efficient representation of sparse connectivity. This is a simple set of bit mapped registers that allow groups of contiguous zeros in memory to be removed from real memory.</p> <p>8. <i>Output Buffer</i>: The Output Buffer contains the arbitration logic for access to the Inter-PN bus and the OutBus. Both 16 and 8 bit values can be transmitted. 16 bit values require 2 clocks to transmit over the 8 bit OutBus. The OutBus arbitration and transmission mechanism operates separately from, but is synchronized with, the SIMD control to allow PNs to both transmit and do multiply-accumulation simultaneously. This capability is required when outputs of a layer are fed back as inputs to the next layer. This feature is called the Virtual PN<sup>5</sup>. There are several OutBus arbitration modes.”</p> <p>Dan Hammerstrom, <i>A VLSI architecture for high-performance, low-cost, on-chip learning</i>, 1990 IJCNN International Joint Conference on Neural Networks, 1990, vol.2, at 539.</p> <p>“Each PN is similar to an individual Digital Signal Processor (DSP). Figure 2 shows a block diagram of a PN in the CNAPS architecture. There are eight functional units in a PN:</p> <ol style="list-style-type: none"> <li>1. The Input unit receives data from the INbus.</li> <li>2. The Output unit sends data to the OUTbus.</li> <li>3. The Multiplier can perform an 8x8 multiply, a 8x16 multiply, or a 16x16 multiply.</li> <li>4. The Adder can perform a 32x32 addition.</li> <li>5. The Logic unit can perform the Boolean logic functions of AND, OR, XOR. The logic unit can also shift words right and left.</li> <li>6. The Register unit contains 32 - 16 bit registers for temporary value storage.</li> <li>7. There is memory local to each PN.</li> <li>8. The Memory Address Unit accesses the local memory.”</li> </ol> <p>Thomas E. Baker, <i>Artificial neural network and image processing using the Adaptive Solutions' architecture</i>, Proc. SPIE 1452, Image Processing Algorithms and Techniques II, June 1, 1991, at 503.</p> <p>“The CNAPS-1064 component is a high-performance processor consisting of 64 individual</p>

## Exhibit 3 - CNAPS

Claim Limitation (Claim 7)	Exemplary Disclosure
	<p>processor nodes (PNs) that operate in parallel. Each PN is a fully functional digital processor with an adder, multiplier, logic unit, local memory, a set of input and output buses, two internal buses that connect functional units, and an inter-PN bus that connects the PNs serially.” Adaptive Solutions, Inc., <i>CNAPS Data Book</i>, CNAPS Hardware Series, October 15, 1993, at HAMMERSTROM-00000158.</p> <p>“The CNAPS PN array is the processing engine for a CNAPS system. The array consists of one or more CNAPS chips, each containing up to 64 processors (PNs). Each PN is a complete fixed-point processor with its own 4KB of on-chip memory. A PN can perform 1-, 8-, or 16-bit integer arithmetic and can execute a multiply-and-accumulate operation in one clock cycle. The effect is similar to having 64 digital signal processors (DSPs) on one chip.” Adaptive Solutions, Inc., <i>Getting Acquainted with CNAPS</i>, CNAPS Core Series, October 15, 1993, at HAMMERSTROM-00000490.</p> <p><u>As it relates to the Court’s construction of LPHDR execution unit, CNAPS included both addressable memory paired with the processing element(s) and control for the processing elements (to the extent that control is interpreted to include any signaling that affects the operation of the processing element). See, e.g., Dan Hammerstrom, AN INTRODUCTION TO NEURAL AND ELECTRONIC NETWORKS 337 (2d ed. 1995) (“As Figure 2 suggests, each PN has a local memory,4 a multiplier, an adder/subtractor, a shifter/logic unit, a register tile,’ and a memory addressing unit....”); id. at 336 (“The CNAPS architecture consists of an array of processors controlled by a sequencer...”); Dan Hammerstrom, A VLSI architecture for high-performance, low-cost, on-chip learning, 1990 IJCNN International Joint Conference on Neural Networks, 1990, vol.2, at 539 (“Each Processor Node has internal units connected via control signals and several data buses...”).</u></p> <p><u>To the extent Singular contends that the CNAPS did not include addressable memory paired with the processing element(s) and control for the processing elements, processing elements that were paired with addressable memory and control for the processing elements (to the extent that control is interpreted to include any signaling that affects the operation of the processing element) were well-known in the art, as explained in Section IV.C.1.d of the Amended Responsive Contentions Regarding Non-Infringement and Invalidity. See, e.g., ’273</u></p>

## Exhibit 3 - CNAPS

Claim Limitation (Claim 7)	Exemplary Disclosure
	<p><u>patent, 3:49-56 (describing admitted prior art). To the extent Singular nonetheless contends that one of skill in the art would have needed a motivation to combine CNAPS with processing elements with paired addressable memory and/or processing elements with control, one of skill in the art would have been motivated to do so based on the teachings of any of Dockser, Belanović, Belanović and Leeser, Shirazi, Lienhart, the admitted prior art, Patterson &amp; Hennessy, in Computer Organization &amp; Design, The Hardware Software Interface (3d. Ed. 2005), and/or Hamada.</u></p>
<p>[156c] wherein the dynamic range of the possible valid inputs to the first operation is at least as wide as from 1/1,000,000 through 1,000,000 and for at least <math>X=5\%</math> of the possible valid inputs to the first operation, the statistical mean, over repeated execution of the first operation on each specific input from the at least <math>X\%</math> of the possible valid inputs to the first operation, of the numerical values represented by the first output signal of the LPHDR unit executing the first operation on that input differs by at least <math>Y=0.05\%</math> from the result of an exact mathematical calculation of the first operation on the numerical values of that same input; and</p>	<p>As explained below and in the Responsive Contentions Regarding Non-Infringement and Invalidity (“Responsive Contentions”), it would have been obvious to one of skill in the art based on the disclosures in CNAPS (alone or in combination with the reduced precision floating point teachings of Dockser, Tong, Belanovic / Belanovic and Leeser, Lee, Shirazi, Aty, Sudha, and TMS320C32, or the logarithmic format disclosed in GRAPE-3 and Hoefflinger) that the dynamic range of the possible valid inputs to the first operation is at least as wide as from 1/1,000,000 through 1,000,000 and for at least <math>X=5\%</math> of the possible valid inputs to the first operation, the statistical mean, over repeated execution of the first operation on each specific input from the at least <math>X\%</math> of the possible valid inputs to the first operation, of the numerical values represented by the first output signal of the LPHDR unit executing the first operation on that input differs by at least <math>Y=0.05\%</math> from the result of an exact mathematical calculation of the first operation on the numerical values of that same input.</p> <p>CNAPS was designed to perform neural network algorithms at the minimum required precision.</p> <p>“Another major simplification of the CNAPS architecture, which is found in other digital ANN chips, was the use of limited-precision, fixed-point arithmetic. Many researchers have shown that floating point and high precision are unnecessary in ANN simulation (Fahlman and Hoehfeld, 1992). CNAPS supported 1-, 8-, and 16-bit precision in hardware. Consequently, the PNs were smaller and cheaper. This reduced precision was more than adequate for the applications implemented on CNAPS.” Dan Hammerstrom, <i>Digital VLSI for Neural Networks</i>, at 9-10, <a href="https://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.110.7541&amp;rep=rep1&amp;type=pdf">https://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.110.7541&amp;rep=rep1&amp;type=pdf</a>.</p>

## Exhibit 3 - CNAPS

Claim Limitation (Claim 7)	Exemplary Disclosure
	<p>Reduced precision computations were commonly used in connection with neurocomputers like the CNAPS. <i>See, e.g.:</i></p> <p>“BP implementations typically use 32-bit floating point math. This largely eliminates scaling, precision and dynamic range issues. Efficient hardware implementation dictates integer arithmetic units with precision no greater than required. Baker [Bak90] has shown 16-bit integer weights are sufficient for BP training and much lower values adequate for use after training.” Hal McCartor, <i>Back Propagation Implementation on the Adaptive Solutions CNAPS Neurocomputer Chip</i>, ADVANCES IN NEURAL INFORMATION PROCESSING SYSTEMS 3, 1990, at 1029.</p> <p>“The inherent parallelism in ANN and connectionist models suggests an opportunity to speed up the simulations. Their simple, low precision computations also suggest an opportunity to employ simpler and cheaper, low-precision digital hardware implemented by full-custom silicon or by FPGAs (Field Programmable Gate Arrays).” Dan Hammerstrom, <i>Digital VLSI for Neural Networks</i>, at 2,  <a href="https://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.110.7541&amp;rep=rep1&amp;type=pdf">https://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.110.7541&amp;rep=rep1&amp;type=pdf</a>.</p> <p>“Conforms to idea that [Artificial Neural Network] consists of large numbers of simple, low precision processing elements. . . NetTalk, for example, showed no appreciable difference between 32 bit floating point and 15 bit integer.” Dan Hammerstrom &amp; Gary Tahara, <i>CNAPS (Connected Network of Adaptive ProcessorS)</i>, Presentation at Hot Chips Conference, July 8, 1991, at 3, <a href="https://www.hotchips.org/wp-content/uploads/hc_archives/hc03/3_Tue/Hc3.S7/Hc3.7.2.pdf">https://www.hotchips.org/wp-content/uploads/hc_archives/hc03/3_Tue/Hc3.S7/Hc3.7.2.pdf</a>.</p> <p>“The limited precision simulation results show that sixteen total bits of precision (one sign bit, three bits to the left of the binary point and twelve bits to the right) are adequate. . . . A floating point representation may be able to learn successfully with fewer than twelve bits in the mantissa, but floating point units are too expensive for massively parallel architectures.” Tom Baker &amp; Dan Hammerstrom, <i>Modifications to Artificial Neural Networks Models for Digital Hardware Implementation</i>, Oregon Graduate Center, Technical Report No. CS/E 88-035, December 1988, <a href="https://scholararchive.ohsu.edu/downloads/pc289j41v?locale=en">https://scholararchive.ohsu.edu/downloads/pc289j41v?locale=en</a>.</p>

**Exhibit 3 - CNAPS**

Claim Limitation (Claim 7)	Exemplary Disclosure
	<p>Accordingly, CNAPS discloses the use of arithmetic units designed to perform calculations using reduced-precision 1-, 8-, and 16-bit math for calculations that typically use 32-bit floating point math, operating on only the 1-, 8-, or 16 most significant bits in the registers. <i>See, e.g.:</i></p> <p>“The X1 has three basic weight modes: 1 bit, 8 bits (7 mantissa + sign), and 16 bits (15 mantissa + sign). In 8 and 16 bit modes, there is a sign bit. When the mantissa is zero, it is generally considered a null weight. Most algorithms require 8 bits, but there are some algorithms that need only single bit precision and some that require up to 16 bits. Bit mode does not have an explicit sign but can be run twice with both positive and negative weight vectors. Two’s complement representation is used throughout. Many of the units recognize and deal with negative numbers appropriately, other units assume unsigned values, that is, the sign bit is ignored.</p> <p>Values transmitted onto and off the chip can be either 8 or 16 bits (16 bit values require two uses of the 8 bit buses). The preferred mode is 8 bits, since for most neural network algorithms that is sufficient. The 16 bit mode is provided for those applications, such as digital signal processing or the loading of 16 bit weights, that need more precision. With 16 bit values, bit 15 is always the sign bit. 8 bit values can be switched to signed or unsigned.” Dan Hammerstrom &amp; Hal McCartor, <i>X1 Programmer’s Guide and Reference Manual</i>, Adaptive Solutions, Inc., March 7, 1990, at HAMMERSTROM-00000276.</p> <p>“Each PN provides four data precisions for representing data: 1-bit, 8-bit (signed or unsigned), 16-bit (signed or unsigned) and 32-bit (signed). Most algorithms use 8 bits, but some algorithms need only single-bit precision and some require up to 16 bits. Single-bit data does not have an explicit sign, but the network can be run twice with both positive and negative data values.” Adaptive Solutions, Inc., <i>CNAPS Data Book</i>, CNAPS Hardware Series, October 15, 1993, at HAMMERSTROM-00000154.</p> <p>CNAPS further discloses expanding the dynamic range of the possible valid inputs to an operation by reducing the precision of the operation’s calculations.</p>

## Exhibit 3 - CNAPS

Claim Limitation (Claim 7)	Exemplary Disclosure
	<p>“In place of floating point variables, the CNAPS provides scaled variables, of one or two bytes, in which a fixed point can be set at a given bit to separate the whole number part from the fraction. Choosing the fixed point position is an important consideration. A low fixed point reduces the precision, whereas a high fixed point to achieve greater precision may result in overflows of the whole part. CNAPS programs often have multiple scaled variables with different fixed points to be used where greater or less fractional precision is necessary.” Jason M. Kinser and Thomas Lindblad, <i>Implementation of pulse-coupled neural networks in a CNAPS environment</i>, IEEE TRANSACTIONS ON NEURAL NETWORKS, vol. 10, no. 3, May 1999.</p> <p>“In the CNAPS real numbers are represented by a scaled type. The scaled value requires the user to define the number of bits that will represent the integer part of the real number and the number of bits that will represent the decimal portion of the number with the total number of bits being 16 or less. Since this is not the mantissa/exponent representation the dynamic range of the scaled type is quite limited.” Jason M. Kinser and Thomas Lindblad, <i>Implementation of pulse-coupled neural networks in a CNAPS environment</i>, IEEE TRANSACTIONS ON NEURAL NETWORKS, vol. 10, no. 3, May 1999.</p> <p>“The CNAPS architecture uses fixed-point arithmetic, so floating-point values must be converted (or <i>scaled</i>) to fixed-point values. CNAPS-C supports fixed-point dated through the scaled type construction, which lets you define a fixed-point type with a specified precision.” Adaptive Solutions, Inc., <i>Getting Acquainted with CNAPS</i>, CNAPS Core Series, October 15, 1993, at HAMMERSTROM-00000497.</p> <p>For the reasons explained in the Responsive Contentions, it would have been obvious to one of skill in the art to have substituted the fixed-point number format used in CNAPS for a floating-point format that met the claimed minimum range and precision requirements, and to have used the reduced-precision floating-point number formats disclosed in Tong, Dockser, Belanovic / Belanovic and Leaser, Shirazi, Sudha, Aty, and TMS 320C32 or the logarithmic format disclosed in GRAPE-3 and Hoefflinger, either alone or in combination. <i>See also</i> Appendix to Responsive Contentions (detailing error rates associated with different mantissa sizes).</p>
[156d] at least one first computing	CNAPS discloses at least one first computing device adapted to control the operation of the at

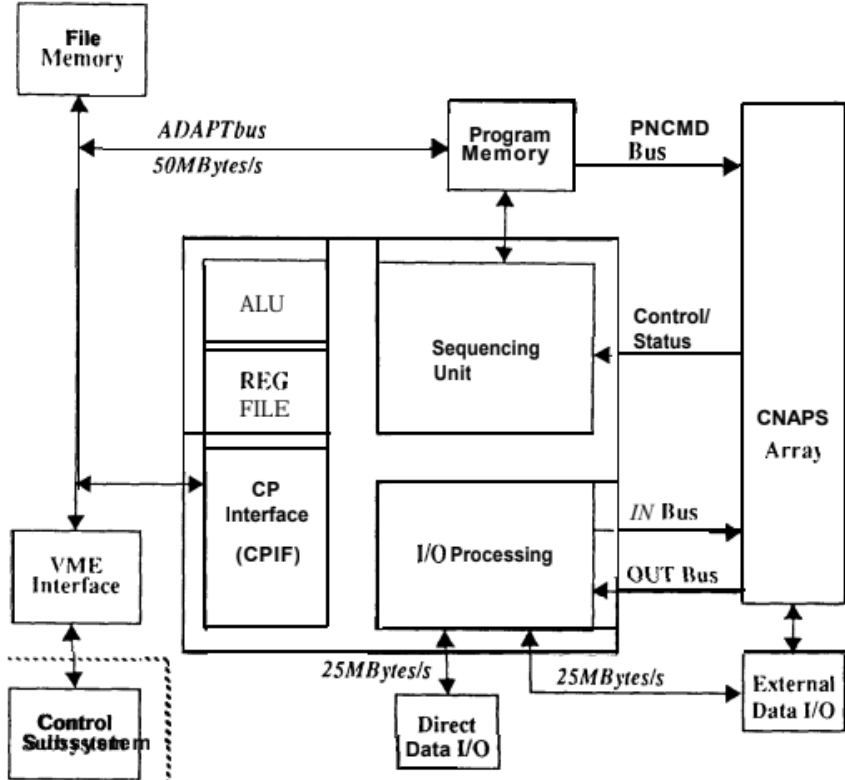
**Exhibit 3 - CNAPS**

Claim Limitation (Claim 7)	Exemplary Disclosure
<p>device adapted to control the operation of the at least one first LPHDR execution unit;</p>	<p>least one first LPHDR execution unit. <i>See, e.g.:</i></p> <p>“The digital CNAPS neurocomputer (Adaptive Solutions, Inc.) has been designed as a general purpose parallel computer for artificial neural networks and image processing tasks and was introduced by Hammerstrom in 1990. CNAPS (Connected Network of Adaptive Processors) consists of a sequencer and a linear array of maximal processor nodes and is a SIMD computer (single instruction stream – multiple data stream). A single instruction bus and three data busses connect the sequencer and the processor nodes: a broadcast input bus (8 bit), a shared output bus (8 bit), and a circular inter-PN bus (2 bit) connecting neighbouring processors (figure 3 left). The sequencer broadcasts an instruction to the processor nodes single instruction stream which synchronously execute it but process the data from their local memories multiple data stream.” Peter Paschke &amp; Ralf Möller, <i>Simulation of Sparse Random Networks on a CNAPS SIMD Neurocomputer</i>, Proc. 15th IMACS World Congress, Berlin, August 1997, at 766.</p> <p>“The Adaptive Solutions CNAPS architecture is embodied in a single chip digital neurocomputer with 64 processors running at 25 megahertz. All processors receive the same instruction which they conditionally execute. Multiplication and addition are performed in parallel allowing 1.6 billion inner product steps per second per chip. Each processor has a 32-bit adder, 9-bit by 16-bit multiplier (16 by 16 in two clock cycles), shifter, logic unit, 32 16-bit registers, and 4096 bytes of local memory. Input and output are accomplished over 8-bit input and output buses common to all processors. The output bus is tied to the input bus so that output of one processor can be broadcast to all others. When multiple chips are used, they appear to the user as one chip with more processors. Special circuits support finding the maximum of values held in each processor and conserving weight space for sparsely connected networks. An accompanying sequencer chip controls instruction flow, input and output.” Hal McCartor, <i>Back Propagation Implementation on the Adaptive Solutions CNAPS Neurocomputer Chip</i>, ADVANCES IN NEURAL INFORMATION PROCESSING SYSTEMS 3, 1990, at 1029.</p> <p>“The CNAPS architecture consists of an array of processors controlled by a sequencer, both implemented as a chip set developed by Adaptive Solutions. Inc. The sequencer is a one-chip device called the CNAPS Sequencer Chip (CSC). The processor array is also a one-chip device,</p>

**Exhibit 3 - CNAPS**

Claim Limitation (Claim 7)	Exemplary Disclosure
	<p>available with either 64 or 16 processors per chip (the CNAPS- 1064 or CNAPS- 10 16).” Dan Hammerstrom, AN INTRODUCTION TO NEURAL AND ELECTRONIC NETWORKS 336 (2d ed. 1995).</p> <p>“CNAPS is a single instruction, multiple data stream (SIMD) architecture. A SIMD computer has one instruction sequencing/control unit and many processor nodes (PNs). In CNAPS, the PNs are connected in a one-dimensional array (Figure 1) in which each PN can “talk” only to its right or left neighbors. The sequencer broadcasts each instruction plus input data to all PNs, which execute the same instruction at each clock. The PNs transmit output data to the sequencer, with several arbitration modes controlling access to the output bus.” Dan Hammerstrom, AN INTRODUCTION TO NEURAL AND ELECTRONIC NETWORKS 337 (2d ed. 1995).</p> <p>“The CSC sequencer (Figure 3) performs program sequencing for the PN array and has private access to a program memory. The CSC also performs input/ output (I/O) processing for the array, writing input data to the array and reading output data from it.” Dan Hammerstrom, AN INTRODUCTION TO NEURAL AND ELECTRONIC NETWORKS 337 (2d ed. 1995).</p> <div data-bbox="898 868 1690 1261" data-label="Diagram"> </div> <p>FIGURE 1 The basic CNAPS architecture. CNAPS is a single instruction, multiple data (SIMD) architecture that uses broadcast input, one-dimensional interprocessor communication, and a single shared output bus.</p> <p>Dan Hammerstrom, AN INTRODUCTION TO NEURAL AND ELECTRONIC NETWORKS 338 (2d ed.</p>

## Exhibit 3 - CNAPS

Claim Limitation (Claim 7)	Exemplary Disclosure
	<p>1995).</p>  <p><b>FIGURE 3</b> The CNAPS sequencer chip (CSC) internal structure. The CSC accesses an external program store, which contains both CSC and CNAPS PN array instructions. PN array instructions are broadcast to all PNs. CSC instructions control sequencing and all array input and output.</p> <p>Dan Hammerstrom, AN INTRODUCTION TO NEURAL AND ELECTRONIC NETWORKS 339 (2d ed. 1995).</p> <p>“A separate sequencer containing an instruction store and microsequencer controls SIMD PN execution. The sequencer places data or literals onto the IN bus and forward sequencing</p>

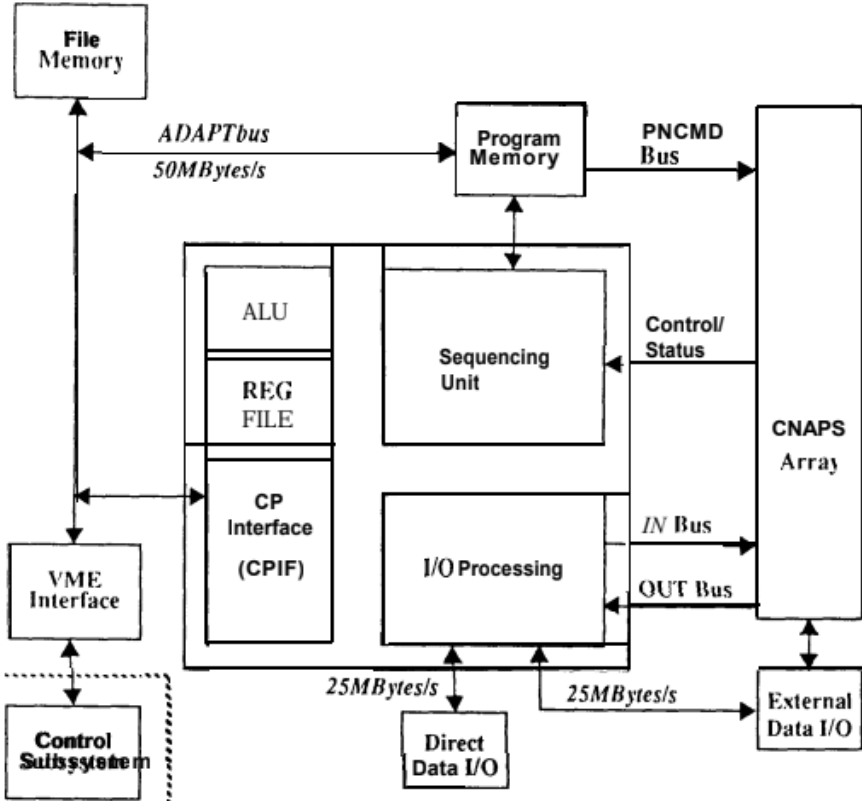
**Exhibit 3 - CNAPS**

Claim Limitation (Claim 7)	Exemplary Disclosure
	<p>commands to the PN array.” Adaptive Solutions, Inc., <i>CNAPS Neurocomputing</i>, 1991, at HAMMERSTROM-00000999.</p>
<p>[156e] wherein the at least one first computing device comprises at least one of a central processing unit (CPU), a graphics processing unit (GPU), a field programmable gate array (FPGA), a microcode-based processor, a hardware sequencer, and a state machine;</p>	<p>CNAPS discloses at least one first computing device comprises at least one of a central processing unit (CPU), a graphics processing unit (GPU), a field programmable gate array (FPGA), a microcode-based processor, a hardware sequencer, and a state machine. <i>See, e.g.:</i></p> <p>“The digital CNAPS neurocomputer (Adaptive Solutions, Inc.) has been designed as a general purpose parallel computer for artificial neural networks and image processing tasks and was introduced by Hammerstrom in 1990. CNAPS (Connected Network of Adaptive Processors) consists of a sequencer and a linear array of maximal processor nodes and is a SIMD computer (single instruction stream – multiple data stream). A single instruction bus and three data busses connect the sequencer and the processor nodes: a broadcast input bus (8 bit), a shared output bus (8 bit), and a circular inter-PN bus (2 bit) connecting neighbouring processors (figure 3 left). The sequencer broadcasts an instruction to the processor nodes single instruction stream which synchronously execute it but process the data from their local memories multiple data stream.” Peter Paschke &amp; Ralf Möller, <i>Simulation of Sparse Random Networks on a CNAPS SIMD Neurocomputer</i>, Proc. 15th IMACS World Congress, Berlin, August 1997, at 766.</p> <p>“The Adaptive Solutions CNAPS architecture is embodied in a single chip digital neurocomputer with 64 processors running at 25 megahertz. All processors receive the same instruction which they conditionally execute. Multiplication and addition are performed in parallel allowing 1.6 billion inner product steps per second per chip. Each processor has a 32-bit adder, 9-bit by 16-bit multiplier (16 by 16 in two clock cycles), shifter, logic unit, 32 16-bit registers, and 4096 bytes of local memory. Input and output are accomplished over 8-bit input and output buses common to all processors. The output bus is tied to the input bus so that output of one processor can be broadcast to all others. When multiple chips are used, they appear to the user as one chip with more processors. Special circuits support finding the maximum of values held in each processor and conserving weight space for sparsely connected networks. An accompanying sequencer chip controls instruction flow, input and output.” Hal McCartor, <i>Back Propagation Implementation on the Adaptive Solutions CNAPS Neurocomputer Chip</i>, ADVANCES IN NEURAL INFORMATION PROCESSING SYSTEMS 3, 1990, at</p>

## Exhibit 3 - CNAPS

Claim Limitation (Claim 7)	Exemplary Disclosure
	<p>1029.</p> <p>“The CNAPS architecture consists of an array of processors controlled by a sequencer, both implemented as a chip set developed by Adaptive Solutions, Inc. The sequencer is a one-chip device called the CNAPS Sequencer Chip (CSC). The processor array is also a one-chip device, available with either 64 or 16 processors per chip (the CNAPS- 1064 or CNAPS- 10 16).” Dan Hammerstrom, AN INTRODUCTION TO NEURAL AND ELECTRONIC NETWORKS 336 (2d ed. 1995).</p> <p>“CNAPS is a single instruction, multiple data stream (SIMD) architecture. A SIMD computer has one instruction sequencing/control unit and many processor nodes (PNs). In CNAPS, the PNs are connected in a one-dimensional array (Figure 1) in which each PN can “talk” only to its right or left neighbors. The sequencer broadcasts each instruction plus input data to all PNs, which execute the same instruction at each clock. The PNs transmit output data to the sequencer, with several arbitration modes controlling access to the output bus.” Dan Hammerstrom, AN INTRODUCTION TO NEURAL AND ELECTRONIC NETWORKS 337 (2d ed. 1995).</p> <div data-bbox="905 868 1696 1258" data-label="Diagram"> </div> <p>FIGURE 1 The basic CNAPS architecture. CNAPS is a single instruction, multiple data (SIMD) architecture that uses broadcast input, one-dimensional interprocessor communication, and a single shared output bus.</p> <p>Dan Hammerstrom, AN INTRODUCTION TO NEURAL AND ELECTRONIC NETWORKS 338 (2d ed.</p>

## Exhibit 3 - CNAPS

Claim Limitation (Claim 7)	Exemplary Disclosure
	<p>1995).</p>  <p><b>FIGURE 3</b> The CNAPS sequencer chip (CSC) internal structure. The CSC accesses an external program store, which contains both CSC and CNAPS PN array instructions. PN array instructions are broadcast to all PNs. CSC instructions control sequencing and all array input and output.</p> <p>Dan Hammerstrom, AN INTRODUCTION TO NEURAL AND ELECTRONIC NETWORKS 339 (2d ed. 1995).</p>
[156f] and, wherein the number of LPHDR execution units in the	CNAPS discloses the number of LPHDR execution units in the device exceeds by at least one hundred the non-negative integer number of execution units in the device adapted to execute at

**Exhibit 3 - CNAPS**

Claim Limitation (Claim 7)	Exemplary Disclosure
<p>device exceeds by at least one hundred the non-negative integer number of execution units in the device adapted to execute at least the operation of multiplication on floating point numbers that are at least 32 bits wide.</p>	<p>least the operation of multiplication. <i>See, e.g.:</i></p> <p>“To increase the performance of the CNAPS array, multiple chips can be linked together using a common sequencer. Signals between PNs that cross chip boundaries are transparent to the programmer. The programmer views a multiple chip system as if it were a monolithic array of PNs (see Figure 1). A system with eight CNAPS chips has 512 PNs, and can perform 10 Billion MAC operations per second when operating at 20 MHz.” Thomas E. Baker, <i>Artificial neural network and image processing using the Adaptive Solutions’ architecture</i>, Proc. SPIE 1452, Image Processing Algorithms and Techniques II, June 1, 1991, at 505.</p> <p>“The CNAPS architecture had 64 [processor nodes] per chip. . . . Multiple chips could be combined to create larger, more powerful systems.” Dan Hammerstrom, <i>Digital VLSI for Neural Networks</i>, at 10, <a href="https://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.110.7541&amp;rep=rep1&amp;type=pdf">https://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.110.7541&amp;rep=rep1&amp;type=pdf</a>.</p> <p>“The Adaptive Solutions CNAPS architecture chip is a general purpose neurocomputer chip. It has 64 processors, each with 4 K bytes of local memory, running at 25 megahertz.” Hal McCartor, <i>Back Propagation Implementation on the Adaptive Solutions CNAPS Neurocomputer Chip</i>, ADVANCES IN NEURAL INFORMATION PROCESSING SYSTEMS 3, 1990, at 1029.</p> <p>“When multiple chips are used, they appear to the user as one chip with more processors.” Hal McCartor, <i>Back Propagation Implementation on the Adaptive Solutions CNAPS Neurocomputer Chip</i>, ADVANCES IN NEURAL INFORMATION PROCESSING SYSTEMS 3, 1990, at 1029.</p> <p>“Multiple chips can be arbitrarily combined to create large systems.” Dan Hammerstrom &amp; Gary Tahara, <i>CNAPS (Connected Network of Adaptive ProcessorS)</i>, Presentation at Hot Chips Conference, July 8, 1991, <a href="https://www.hotchips.org/wp-content/uploads/hc_archives/hc03/3_Tue/HC3.S7/HC3.7.2.pdf">https://www.hotchips.org/wp-content/uploads/hc_archives/hc03/3_Tue/HC3.S7/HC3.7.2.pdf</a>.</p> <p>“The processor array is also a one-chip device, available with either 64 or 16 processors per</p>

## Exhibit 3 - CNAPS

Claim Limitation (Claim 7)	Exemplary Disclosure
	<p>chip (the CNAPS- 1064 or CNAPS- 10 16). The CSC can control up to eight 1064s or 1016s, which act like one large device.” Dan Hammerstrom, AN INTRODUCTION TO NEURAL AND ELECTRONIC NETWORKS 336 (2d ed. 1995).</p> <p>“Chip boundaries across the array are arbitrary, since the user only sees a single array of [processor nodes]. Thus [CNAPS] systems are easily scalable.” Dan Hammerstrom, <i>A VLSI architecture for high-performance, low-cost, on-chip learning</i>, 1990 IJCNN International Joint Conference on Neural Networks, 1990, vol.2, at 539.</p> <p>“The CNAPS-1064 chip is a CMOS device with the parallel processing power of a DSP and the flexibility to program it for a wide variety of applications. Each CNAPS-1064 component contains 64 interconnected processors, called processor nodes or <i>PNs</i>, which can be programmed to operate in parallel on a given set of data, Each PN is a simple arithmetic processor with its own local memory, adder, multiplier, logic unit, register file, and I/O buffers.</p> <p>Multiple CNAPS-1064 chips can be strung together to dramatically improve the speed and processing power of the system, Such a connection of CNAPS-1064 chips is called the <i>CNAPS array</i> (or <i>PN array</i>). PNs in the array are numbered sequentially starting with zero (0) and continue through to the last PN on the last CNAPS chip, The PN array is controlled by a sequencing and I/O processing unit the CNAPS CSC (See Part I for details on the CSC.)” Adaptive Solutions, Inc., <i>CNAPS Data Book</i>, CNAPS Hardware Series, October 15, 1993, at HAMMERSTROM-00000153.</p> <p>“The CNAPS system is a complete parallel computer system based on a custom digital chip containing up to 64 processors. Current CNAPS systems have one to eight CNAPS chips for a total of 64 to 512 processors. Each processor can execute a multiple-and-accumulate (dot product) operation in one clock cycle, so an array of 512 processors yields up to 10.24 <i>billion</i> multiply-accumulates per second at 20 MHz. The result is similar to having up to 512 digital signal processors (DSPs) in one system. CNAPS provides the computational power needed for real-world, real-time pattern recognition problems—all in a compact, economical package.” Adaptive Solutions, Inc., <i>Getting Acquainted with CNAPS</i>, CNAPS Core Series, October 15, 1993, at HAMMERSTROM-00000482.</p>

**Exhibit 3 - CNAPS**

Claim Limitation (Claim 7)	Exemplary Disclosure
	<p>“The CNAPS hardware consists of a linear array of single-instruction, multiple-data processors that operate in parallel on a data set. Each processor operates on its own slice of the data. The CNAPS-C language supports this parallelism.</p> <p>The number of processors is conceptually unlimited: in practice, the number of CNAPS chips determines the number of processors.” Adaptive Solutions, Inc., <i>CNAPS-C Programming Guide</i>, CNAPS Series, March 15, 1996, at HAMMERSTROM-00000567-568.</p> <p>“The CNAPS Server has 256 processors (PNs; processing nodes) that are used to simulate neural network nodes.” Adaptive Solutions, Inc., <i>CNAPS Neurocomputing</i>, 1991 at HAMMERSTROM-00001003.</p>

**Exhibit 3 - CNAPS****'273 Patent**

<b>Claim Limitation (Claim 53)</b>	<b>Exemplary Disclosure</b>
<b>[273a]</b> A device:	CNAPS discloses a device. Specifically, the CNAPS is computer architecture consisting of an array of processors controlled by a sequencer. <i>See [156a]</i> .
<b>[273b]</b> comprising at least one first low precision high-dynamic range (LPHDR) execution unit adapted to execute a first operation on a first input signal representing a first numerical value to produce a first output signal representing a second numerical value,	CNAPS discloses at least one first low precision high dynamic range (LPHDR) execution unit adapted to execute a first input signal representing a first numerical value to produce a first output signal representing a second numerical value. <i>See [156b]</i> .
<b>[273c]</b> wherein the dynamic range of the possible valid inputs to the first operation is at least as wide as from 1/1,000,000 through 1,000,000 and for at least X=5% of the possible valid inputs to the first operation, the statistical mean, over repeated execution of the first operation on each specific input from the at least X % of the possible valid inputs to the first operation, of the numerical values represented by the first output signal of the LPHDR unit executing the first operation on that input differs by at least Y=0.05% from the result of an exact mathematical calculation of the first operation on the numerical values of that same input;	For the reasons explained above and in the Responsive Contentions, it would have been obvious to one of skill in the art based on the disclosures in CNAPS (alone or in combination with the reduced precision floating point teachings of Dockser, Tong, Belanovic / Belanovic and Leaser, Lee, Shirazi, Aty, Sudha, and TMS320C32 or the logarithmic format disclosed in GRAPE-3 and Hoefflinger) that the dynamic range of the possible valid inputs to the first operation is at least as wide as from 1/1,000,000 through 1,000,000 and for at least X=5% of the possible valid inputs to the first operation, the statistical mean, over repeated execution of the first operation on each specific input from the at least X% of the possible valid inputs to the first operation, of the numerical values represented by the first output signal of the LPHDR unit executing the first operation on that input differs by at least Y=0.05% from the result of an exact mathematical calculation of the first operation on the numerical values of that same input. <i>See [156c]</i> .
<b>[273d]</b> wherein the number of LPHDR execution units in the device exceeds by at least one hundred the non-negative integer number of execution units in the device adapted to execute at least the operation of multiplication on floating point numbers that are at least 32 bits wide.	CNAPS discloses the number of LPHDR execution units in the device exceeds by at least one hundred the non-negative integer number of execution units in the device adapted to execute at least the operation of multiplication on floating point numbers that are at least 32 bits wide. <i>See [156f]</i> .

**Exhibit 3 - CNAPS**

## '961 Patent

Claim Limitation (Claim 4)	Exemplary Disclosure
<b>[961a]</b> A device comprising:	CNAPS discloses a device. Specifically, the CNAPS is computer architecture consisting of an array of processors controlled by a sequencer. <i>See</i> [156a].
<b>[961b]</b> at least one first low precision high-dynamic range (LPHDR) execution unit adapted to execute a first operation on a first input signal representing a first numerical value to produce a first output signal representing a second numerical value,	CNAPS discloses at least one first low precision high dynamic range (LPHDR) execution unit adapted to execute a first input signal representing a first numerical value to produce a first output signal representing a second numerical value. <i>See</i> [156b].
<b>[961c]</b> wherein the dynamic range of the possible valid inputs to the first operation is at least as wide as from 1/1,000,000 through 1,000,000 and for at least X=10% of the possible valid inputs to the first operation, the statistical mean, over repeated execution of the first operation on each specific input from the at least X% of the possible valid inputs to the first operation, of the numerical values represented by the first output signal of the LPHDR unit executing the first operation on that input differs by at least Y=0.2% from the result of an exact mathematical calculation of the first operation on the numerical values of that same input; and	For the reasons explained above and in the Responsive Contentions, it would have been obvious to one of skill in the art based on the disclosures in CNAPS (alone or in combination with the reduced precision floating point teachings of Dockser, Tong, Belanovic / Belanovic and Leaser, Lee, Shirazi, Aty, Sudha, and TMS320C32 or the logarithmic format disclosed in GRAPE-3 and Hoefflinger) that the dynamic range of the possible valid inputs to the first operation is at least as wide as from 1/1,000,000 through 1,000,000 and for at least X=10% of the possible valid inputs to the first operation, the statistical mean, over repeated execution of the first operation on each specific input from the at least X% of the possible valid inputs to the first operation, of the numerical values represented by the first output signal of the LPHDR unit executing the first operation on that input differs by at least Y=0.2% from the result of an exact mathematical calculation of the first operation on the numerical values of that same input. <i>See</i> [156c].
<b>[961d]</b> at least one first computing device adapted to control the operation of the at least one first LPHDR execution unit.	CNAPS discloses at least one first computing device adapted to control the operation of the at least one first LPHDR execution unit. <i>See</i> [156d].

**Exhibit 3 - CNAPS**

Claim Limitation (Claim 13)	Exemplary Disclosure
[961e] A device comprising:	CNAPS discloses a device. Specifically, the CNAPS is computer architecture consisting of an array of processors controlled by a sequencer. <i>See</i> [156a].
[961f] a plurality of components comprising:	<p>CNAPS discloses at least one first low precision high dynamic range (LPHDR) execution unit adapted to execute a first input signal representing a first numerical value to produce a first output signal representing a second numerical value. <i>See</i> [156b].</p> <p>CNAPS discloses at least one first computing device adapted to control the operation of the at least one first LPHDR execution unit. <i>See above</i> [156d].</p>
[961g] at least one first low precision high-dynamic range (LPHDR) execution unit adapted to execute a first operation on a first input signal representing a first numerical value to produce a first output signal representing a second numerical value,	CNAPS discloses at least one first low precision high dynamic range (LPHDR) execution unit adapted to execute a first input signal representing a first numerical value to produce a first output signal representing a second numerical value. <i>See</i> [156b].
[961h] wherein the dynamic range of the possible valid inputs to the first operation is at least as wide as from 1/1,000,000 through 1,000,000 and for at least X=10% of the possible valid inputs to the first operation, the statistical mean, over repeated execution of the first operation on each specific input from the at least X% of the possible valid inputs to the first operation, of the numerical values represented by the first output signal of the LPHDR unit executing the first operation on that input differs by at least Y=0.2% from the result of an exact mathematical calculation of the first operation on the numerical values of that same input.	For the reasons explained above and in the Responsive Contentions, it would have been obvious to one of skill in the art based on the disclosures in CNAPS (alone or in combination with the reduced precision floating point teachings of Dockser, Tong, Belanovic / Belanovic and Leeser, Lee, Shirazi, Aty, Sudha, and TMS320C32 or the logarithmic format disclosed in GRAPE-3 and Hoefflinger) that the dynamic range of the possible valid inputs to the first operation is at least as wide as from 1/1,000,000 through 1,000,000 and for at least X=10% of the possible valid inputs to the first operation, the statistical mean, over repeated execution of the first operation on each specific input from the at least X% of the possible valid inputs to the first operation, of the numerical values represented by the first output signal of the LPHDR unit executing the first operation on that input differs by at least Y=0.2% from the result of an exact mathematical calculation of the first operation on the numerical values of that same input. <i>See</i> [156c].